

## 1 Getting Started

**Read through this page carefully.** You may typeset your homework in latex or submit neatly handwritten/scanned solutions. Please start each question on a new page. Deliverables:

1. Submit a PDF of your writeup, **with an appendix for your code**, to assignment on Gradescope, “HW13 Write-Up”. If there are graphs, include those graphs in the correct sections. Do not simply reference your appendix.
2. If there is code, submit all code needed to reproduce your results, “HW13 Code”.
3. If there is a test set, submit your test set evaluation results, “HW13 Test Set”.

After you’ve submitted your homework, watch out for the self-grade form.

- (a) Who else did you work with on this homework? In case of course events, just describe the group. How did you work on this homework? Any comments about the homework?

- (b) Please copy the following statement and sign next to it. We just want to make it *extra* clear so that no one inadvertently cheats.

*I certify that all solutions are entirely in my words and that I have not looked at another student’s solutions. I have credited all external sources in this write up.*

## 2 Gradient boosting and early stopping

In this problem we show how the greedy algorithms Matching Pursuit (a simplified version of the Orthogonal Matching Pursuit that you might have seen in EE16A) and AdaBoost which you have seen in class, can both be interpreted as taking steps greedily in a direction which is closest to the negative gradient in a restricted space. We start by connecting gradient descent on the square loss with Matching Pursuit.

**Gradient descent on square loss** So far in lecture, we have considered the squared error loss  $L(\mathbf{w}) := \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$  as a function of the variable  $\mathbf{w}$ . When running gradient descent we took the derivative of  $L$  with respect to  $\mathbf{w}$  and obtained iterations

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \alpha_t \mathbf{X}^\top (\mathbf{X}\mathbf{w}^t - \mathbf{y}) \quad (1)$$

We can substitute  $\mathbf{v} = \mathbf{X}\mathbf{w} \in \mathbb{R}^n$  and equivalently consider the loss function as a map  $\mathcal{L} : \mathcal{V} \mapsto \mathbb{R}$  mapping from the Euclidean space  $\mathcal{V} = \mathbb{R}^n$  to real values:

$$\mathcal{L}(\mathbf{v}) = \frac{1}{2} \|\mathbf{y} - \mathbf{v}\|_2^2.$$

In order to minimize this new loss  $\mathcal{L}$  with respect to  $\mathbf{v}$ , we could imagine simply running gradient descent on this loss as follows:

$$\mathbf{v}^{t+1} = \mathbf{v}^t - \alpha_t \nabla_{\mathbf{v}} \mathcal{L}(\mathbf{v})|_{\mathbf{v}=\mathbf{v}^t}, \quad (2)$$

for some suitably chosen step size  $\alpha_t > 0$ .

(a) **Compute the gradient  $\nabla_{\mathbf{v}} \mathcal{L}(\mathbf{v})$ . Write down the gradient descent update for  $\mathbf{v}^{t+1}$ .**

**What is the gradient  $\nabla_{\mathbf{v}} \mathcal{L}(\mathbf{v})$  evaluated at  $\mathbf{v} = \mathbf{X}\mathbf{w}$  which we denote as  $\nabla_{\mathbf{v}} \mathcal{L}(\mathbf{v})|_{\mathbf{v}=\mathbf{X}\mathbf{w}}$ ? For  $\mathbf{v}^t = \mathbf{X}\mathbf{w}^t$ , compute the expression for  $\mathbf{v}^{t+1}$  using update (2).**

Now we compare this update with the one obtained by running the gradient descent update of  $\mathbf{w}$  on the loss  $L$ . In particular, using equation (1), **write down the expression for  $\mathbf{X}\mathbf{w}^{t+1}$ . How do you interpret the difference between the updates  $\mathbf{X}\mathbf{w}^{t+1}$  and  $\mathbf{v}^{t+1}$ ?** Think in terms of underlying subspaces in  $\mathbb{R}^d$  in which the updates lie.

(b) As we see in the previous part, the updates  $\mathbf{X}\mathbf{w}^{t+1}$  and  $\mathbf{v}^{t+1}$  derived by taking gradients of  $L$  and  $\mathcal{L}$  respectively, are not equivalent. In general, we may not be able to run the gradient updates (2) exactly, as there might be constraints on the update we want to make. Sometimes, these constraints are consequences of the problem formulation (in the previous part for example, we wanted a linear fit for  $y$  depending on  $\mathbf{x}$ ). Sometimes, there are further constraints that come from computational considerations, as in the follow-up parts dealing with AdaBoost). In this problem, we will see that if instead of (2) we take a direction closest to the gradient in some restricted space, both procedures actually match.

In the specific example of linear regression, we want to find a vector  $\tilde{\mathbf{v}}^t$  which can be expressed as a linear function of our data matrix  $\mathbf{X}$ . For this purpose we restrict our search space to the

column space of  $\mathbf{X}$  (a subspace of  $\mathbb{R}^n$ ) in which all vectors  $\tilde{\mathbf{v}}$  can be expressed by  $\tilde{\mathbf{v}} = \mathbf{X}\mathbf{w}$  for some vector  $\mathbf{w}$  which has norm smaller than one (we'll see later why this is necessary), i.e. the set

$$\tilde{\mathcal{V}} = \{\tilde{\mathbf{v}} = \mathbf{X}\mathbf{w} : \|\mathbf{w}\|_2 \leq 1, \mathbf{w} \in \mathbb{R}^d\} \subset \mathcal{V}.$$

Notice that the gradient of  $\mathcal{L}$  with respect to  $\mathbf{v}^t$  need not be in this smaller subspace, i.e.  $\nabla_{\mathbf{v}}\mathcal{L}(\mathbf{v}^t) \notin \tilde{\mathcal{V}}$ . In that case, we consider a greedy strategy: We first find the direction in  $\tilde{\mathcal{V}}$  which best aligns with the negative gradient. That is at each iteration, we find

$$\tilde{\mathbf{v}}^t = \arg \max_{\tilde{\mathbf{v}} \in \tilde{\mathcal{V}}} \langle -\nabla\mathcal{L}(\mathbf{v}^t), \tilde{\mathbf{v}} \rangle. \quad (3)$$

and then compute  $\mathbf{v}^{t+1}$  using the update equation

$$\mathbf{v}^{t+1} = \mathbf{v}^t + \alpha_t \tilde{\mathbf{v}}^t, \quad (4)$$

for an appropriate choice of step size  $\alpha_t$ , where  $\tilde{\mathbf{v}}^t$  now belongs to the subset  $\tilde{\mathcal{V}}$ . We now show how this new procedure yields updates in the same direction as performing the gradient step with respect to  $\mathbf{w}$  explicitly.

Again assume that the current iterate is given by  $\mathbf{v}^t = \mathbf{X}\mathbf{w}^t$  for some vector  $\mathbf{w}^t$ . Now, using equation (3) and the expression for  $\nabla_{\mathbf{v}}\mathcal{L}(\mathbf{v})|_{\mathbf{v}=\mathbf{X}\mathbf{w}^t}$  derived previously, **show that  $\tilde{\mathbf{v}}^t$  in (3) is exactly aligned with the vector  $-\mathbf{X}\mathbf{X}^\top(\mathbf{X}\mathbf{w}^t - \mathbf{y})$  (i.e. they are the same up to some constant scalar). Now write the update equation (4) in terms of  $\mathbf{X}$ ,  $\mathbf{w}^t$ ,  $\mathbf{y}$  and  $\alpha_t$ .** Thus, we have seen that gradient descent on the losses  $L$  and  $\mathcal{L}$  are essentially the same (and in fact the same if adjusting the stepsize  $\alpha_t$  accordingly) by choosing the set  $\tilde{\mathcal{V}}$  accordingly.

Hint: For any function  $h$  mapping to a scalar you may use that

$$\arg \max_{\mathbf{v}=\mathbf{X}\mathbf{w}:\|\mathbf{w}\|_2 \leq 1} h(\mathbf{v}) = \mathbf{X} \arg \max_{\|\mathbf{w}\|_2 \leq 1} h(\mathbf{X}\mathbf{w}).$$

- (c) *Towards Matching Pursuit in the gradient descent framework* In this problem we are going to interpret the direction-selection part of Matching Pursuit as a general gradient-type method as described in (b) for an appropriate choice of the set  $\tilde{\mathcal{V}}$ .

As you have learned in lecture and discussion, Orthogonal Matching Pursuit is a greedy method to pick relevant coordinates that can explain the data best. It searches for a direction (among the columns) such that the residual error can be explained best and updates the entire fit in the augmented space. Matching Pursuit (MP) is a simpler variant where the entire fit is not updated at every step, just the fit along the new coordinate that was chosen.

Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  be the training feature matrix,  $\mathbf{y}$  the observed training target variables from which we want to infer (potentially sparse) weights  $\mathbf{w} \in \mathbb{R}^d$ . At every iteration  $t$ , MP picks the column of  $\mathbf{X}$  with maximal absolute inner product with the current residual  $\mathbf{r}^t = \mathbf{y} - \mathbf{X}\mathbf{w}^t$ , i.e.

$$i_t = \arg \max_{i=1,\dots,d} |\langle \mathbf{r}^t, \phi_i(\mathbf{X}) \rangle|. \quad (5)$$

where  $\mathbf{X} \in \mathbb{R}^{n \times d}$  is the data matrix with columns  $\phi_i(\mathbf{X})$  and rows  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , i.e.  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top = (\phi_1(\mathbf{X}), \dots, \phi_d(\mathbf{X}))$ . Note that  $\mathbf{x}_i$  denotes the  $i$ -th data point and  $\phi_i(\mathbf{X})$  denotes the vector of the  $i$ -th features of all the data points (i.e. the  $i$ -th column of the data matrix  $\mathbf{X}$ ).

The updates of MP at each iteration  $t$  then read

$$\mathbf{X}\mathbf{w}^{t+1} = \mathbf{X}\mathbf{w}^t + \frac{\langle \mathbf{r}^t, \phi_{i_t}(\mathbf{X}) \rangle}{\|\phi_{i_t}(\mathbf{X})\|^2} \phi_{i_t}(\mathbf{X}). \quad (6)$$

Now your task in this part is to establish the correspondence between the MP selection rule (5) and the general approach in (restricted) gradient descent given by (3).

In particular, set  $\mathbf{v}^t = \mathbf{X}\mathbf{w}^t$  and **find the set  $\tilde{\mathcal{V}}$  such that the selection represented by (5) can be understood as just being (3) in action, where the optimal direction reads  $\tilde{\mathbf{v}}^t = \text{sign}(\langle \mathbf{r}^t, \phi_{i_t}(\mathbf{X}) \rangle) \phi_{i_t}(\mathbf{X})$ .**

Hint: don't forget how to deal with the absolute value.

- (d) In traditional gradient descent, the rule for choosing the step-size  $\alpha_t$  is left open. One way of doing this is adaptively by *linesearch* — picking the step-size  $\alpha$  that reduces the loss the most. Formally, choosing stepsize  $\alpha_t$  via linesearch with respect to the search direction  $\tilde{\mathbf{v}}^t$  at each iteration for a loss  $\mathcal{L}$  means finding the best  $\alpha$  such that  $\alpha_t = \arg \min_{\alpha \geq 0} \mathcal{L}(\mathbf{v}^t + \alpha \tilde{\mathbf{v}}^t)$ . **Argue that the update (4) with  $\tilde{\mathbf{v}}^t$  computed using equation (3) with the set  $\tilde{\mathcal{V}}$  you found in (c), and obtaining  $\alpha_t$  via linesearch along that  $\tilde{\mathbf{v}}^t$  is the same as doing the MP update (6).**
- (e) (BONUS) *MP as gradient boosting* The two previous greedy principles of matching-pursuit (column selection and linesearch) can be extended to more general settings of non-linear and non-parametric function classes  $\mathcal{F}$  and general loss functions  $\mathcal{L}$  (which need not be squared loss). Consider the general learning problem where again we are given a bunch of pairwise training samples in the form of  $(\mathbf{x}_i, y_i)$  for  $i = 1, \dots, n$ , with  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \mathcal{Y} \subset \mathbb{R}$  and we want to learn a function  $f : \mathbb{R}^d \rightarrow \mathcal{Y}$  from a function space  $\mathcal{F}$  (with possibly non-linear and non-parametric basis elements, such as trees) which we can then use to predict  $y_{test}$  by  $f(\mathbf{x}_{test})$  for a test sample  $\mathbf{x}_{test}$ .

For this purpose we minimize the following (average) loss  $\mathcal{L} : \mathbb{R}^n \rightarrow \mathbb{R}$

$$\mathcal{L}(\mathbf{f}(\{\mathbf{x}\}_{i=1}^n)) = \mathcal{L}(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i)).$$

for some point-wise loss function  $\ell : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}$  for one sample. Notice that the loss maps a vector to a scalar, as in the linear settings. Here, we use  $\mathbf{f}(\{\mathbf{x}\}_{i=1}^n) := (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$  to denote the vector of function values of  $f \in \mathcal{F}$  at the points  $\mathbf{x}_1, \dots, \mathbf{x}_n$ .

As in gradient descent (4), we want to minimize the loss using an iterative algorithm with updates of the form

$$f^{t+1} = f^t + \alpha_t \tilde{f}^t.$$

Note that our notation follows the convention that  $f^t \in \mathcal{F}$  is the current iterate and  $\tilde{f}^t \in \mathcal{W}$  is the direction that we are adding at each step. The main question is how to choose  $\tilde{f}^t$ . Let's

assume that it is for example computationally much simpler to search over a smaller set of functions  $\mathcal{W} \subset \mathcal{F}$  at each timestep.

Similar to the parametric update in (3), we then choose the function  $\tilde{f} \in \mathcal{W}$  for which the function value vector  $\tilde{\mathbf{f}}(\{\mathbf{x}\}_{i=1}^n)$  has the maximal inner product with the negative gradient, i.e.

$$\tilde{f}^t = \arg \max_{\tilde{f} \in \mathcal{W}} \langle -\nabla \mathcal{L}(\mathbf{f}^t(\{\mathbf{x}\}_{i=1}^n)), \tilde{\mathbf{f}}(\{\mathbf{x}\}_{i=1}^n) \rangle \quad (7)$$

where the gradient  $\nabla \mathcal{L}(\mathbf{f}^t(\{\mathbf{x}\}_{i=1}^n))$ ,  $\tilde{\mathbf{f}}(\{\mathbf{x}\}_{i=1}^n) = \nabla_{\mathbf{v}} \mathcal{L}(\mathbf{v})|_{\mathbf{v}=\mathbf{f}^t(\{\mathbf{x}\}_{i=1}^n)}$  which you have already computed above for the square loss.

Now we want to see how Matching Pursuit can be viewed as an instance of this framework, also called *gradient boosting* or *functional gradient descent*. Specifically, **show that the particular choices of**

$$\mathcal{F} = \{f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} : \mathbf{w} \in \mathbb{R}^d\} \text{ and } \mathcal{W} = \{\tilde{f}(\mathbf{x}) = \pm \mathbf{e}_i^\top \mathbf{x} : i = 1, \dots, d\},$$

**where  $\mathbf{e}_i$  is the  $i$ -th standard basis element with a 1 in the  $i$ -th position and zeros elsewhere, yield Matching Pursuit as a form of gradient boosting.** Notice that  $\mathcal{F}$  is the function space of all linear functions and  $\mathcal{W}$  is the set of functions on vectors  $\mathbf{x} \in \mathbb{R}^d$  which pick out signed coordinates of  $\mathbf{x}$ .

Hint: For matching pursuit, what are  $\tilde{\mathbf{f}}^t(\{\mathbf{x}\}_{i=1}^n)$  and  $\mathbf{f}^t(\{\mathbf{x}\}_{i=1}^n)$ ?

- (f) (BONUS) *AdaBoost as gradient boosting* Now let's have a look at non-linear and non-parametric functions such as decision trees for classification. In this problem we will establish that AdaBoost is also an instance of gradient boosting on a particular loss function (which is not the zero-one loss!).

Notice how each decision tree can be viewed as a function  $\tilde{f} : \mathbb{R}^d \rightarrow \{-1, +1\}$ . Since in each iteration we fit a tree to the weighted samples, our function set  $\mathcal{W}$  now corresponds to the set of decision trees with a fixed set of parameters.

Remember the AdaBoost algorithm from lecture and the last homework. At each iteration  $t$ , it upweights each sample  $i$  by a different  $w_i^t$  in every iteration and finds a new tree which solves the reweighted classification problem

$$\tilde{f}_{Ada}^t = \arg \min_{\tilde{f} \in \mathcal{W}} \sum_{i=1}^n w_i^t \mathbb{I}(y_i \neq \tilde{f}(\mathbf{x}_i)) \text{ where } w_i^t = \frac{\exp(-y_i f^t(\mathbf{x}_i))}{\sum_{i=1}^n \exp(-y_i f^t(\mathbf{x}_i))}$$

where  $f^t$  is a linear combination of functions in  $\mathcal{W}$  found by previous iterations and  $\mathbb{I}(A) = 1$  if  $A$  is true and 0 otherwise (the indicator function). Notice that the samples which the current classifier  $f^t$  gets wrong are assigned a larger weight.

The update for the overall classifier reads

$$f^{t+1} = f^t + \alpha_t \tilde{f}_{Ada}^t \quad (8)$$

The classification rule is then obtained by taking the sign of the final iterate at time  $\tau$ , i.e.  $f_{\text{final}}(\mathbf{x}) = \text{sign}(f^\tau(\mathbf{x}))$  which lies in a bigger space  $\mathcal{F} \supset \mathcal{W}$ .

**Show that  $\tilde{\mathbf{f}}_{\text{Ada}}^t(\{\mathbf{x}\}_{i=1}^n) = \tilde{\mathbf{f}}^t(\{\mathbf{x}\}_{i=1}^n)$  in the gradient boosting framework (7) for the loss  $\ell(y, f(\mathbf{x})) = e^{-yf(\mathbf{x})}$ .** This implies that finding the best fit for the reweighted samples is equivalent to fitting the gradient. For the gradient in (7), notice that we again view  $\mathcal{L}$  as a function of a vector  $\mathbf{v}$  and the gradient  $\nabla \mathcal{L}(\tilde{\mathbf{f}}^t(\{\mathbf{x}\}_{i=1}^n)) = \nabla_{\mathbf{v}} \mathcal{L}(\mathbf{v})|_{\mathbf{v}=\tilde{\mathbf{f}}^t(\{\mathbf{x}\}_{i=1}^n)}$ .

Hint: Now note that for  $y \in \{-1, +1\}$  and  $\tilde{f} : \mathbb{R}^d \rightarrow \{-1, +1\}$ , we have

$$y_i \tilde{f}(\mathbf{x}_i) = 2(1 - \mathbb{I}(y_i \neq \tilde{f}(\mathbf{x}_i))) - 1. \quad (9)$$

Hint: Recall that the maximizer does not change if you scale by or add to the function a constant factor which is independent of the object you are searching over, i.e.

$$\arg \max_{\tilde{f} \in \mathcal{W}} CH(\tilde{\mathbf{f}}(\{\mathbf{x}\}_{i=1}^n)) + G(\mathbf{f}(\{\mathbf{x}\}_{i=1}^n)) = \arg \max_{\tilde{f} \in \mathcal{W}} H(\tilde{\mathbf{f}}(\{\mathbf{x}\}_{i=1}^n))$$

for some functions  $H, G$  that depend on the function value vectors  $\mathbf{f}, \tilde{\mathbf{f}}$ , and a scalar  $C \in \mathbb{R}$ .  $\mathbf{f}$  does not depend on  $\tilde{f}$ .

- (g) (BONUS) *AdaBoost weight as linesearch stepsize* In this problem we show that the weights used in the AdaBoost algorithm at each time step  $t$  correspond to the “best” stepsize in the direction  $\tilde{f}_{\text{Ada}}^t$  using linesearch.

Let’s define the error of the classifier  $\tilde{f}_{\text{Ada}}^t$  as  $\epsilon_t = \sum_{i=1}^n w_i^t \mathbb{I}(y \neq \tilde{f}_{\text{Ada}}^t(\mathbf{x}_i))$ . In practice, AdaBoost uses the stepsize  $\alpha_t = \frac{1}{2} \ln \frac{1-\epsilon_t}{\epsilon_t}$ .

**Show that the AdaBoost stepsize at any time  $t$  corresponds to the stepsize chosen via linesearch for the direction  $\tilde{\mathbf{v}} = \tilde{\mathbf{f}}_{\text{Ada}}^t(\{\mathbf{x}\}_{i=1}^n)$ .**

- (h) Now we will explore how gradient boosting does in practice for both examples we have theoretically studied in the previous sections: Matching Pursuit and AdaBoost. In particular, we will see how the number of updates, equivalent to the number of trees that we are fitting, effects the test and training error.

We have written some code (see the file: `adaboost.py`) that trains several classifiers on the Spam dataset you worked with last week. **Run the code for this part of the question.** It generates a training error plot that uses several different decision trees (depths 1, 2, and 4) as well as several AdaBoost-trained classifiers, each built off of one of these decision tree classifiers. **Do you see a trend in the performance of different trees by themselves? Do you observe a trend in the training error as you use deeper trees for AdaBoost? Why might this happen?**

- (i) Now we examine the test error, and compare against a baseline classifier (a decision tree of depth 9). **Run the code for this part. Is there a difference in the training and test error? Which decision tree depth works best for AdaBoost? Explain your observations.**

- (j) In the last part, you noticed that the test error decreases as a function of boosting iterations in the beginning but eventually it starts to increase when the number of decision trees in Adaboost is pretty large. In practice, this phenomenon motivates us to stop training early and limit the number of classifiers used in a boosting setting. Here “stopping early” means that training error has not reduced to zero but we have stopped training! Refer to the plot from the previous part and answer: **Do you think limiting the number of base classifiers used for AdaBoost would help? Which base classifier can we run more boosting iterations on before the test error starts increasing? Justify your answer intuitively.**
- (k) In this part, we connect this phenomenon to matching pursuit. The provided code generates a synthetic dataset:

$$\begin{aligned} \mathbf{y}_{\text{train}}, \mathbf{y}_{\text{test}} &\in \mathbb{R}^n \\ \mathbf{X}_{\text{train}}, \mathbf{X}_{\text{test}} &\in \mathbb{R}^{n \times d} \\ \mathbf{w} &\in \mathbb{R}^d, \end{aligned}$$

with (for both test and train data)

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \mathbf{z}.$$

Here,  $\mathbf{w}$  is a sparse vector (it has only a few non-zero entries) and  $\mathbf{z}$  denotes noise. Given the observations  $\mathbf{y}_{\text{train}}$  and the feature matrix  $\mathbf{X}_{\text{train}}$ , we are tasked with finding a sparse solution  $\hat{\mathbf{w}}$  for which we use the matching pursuit algorithm.

Let  $\hat{\mathbf{w}}_{\text{MP}}^t$  denote the estimate of  $\mathbf{w}$  recovered by MP after  $t$ -iterations. The code also plots the training error  $\|\mathbf{y}_{\text{train}} - \mathbf{X}_{\text{train}}\hat{\mathbf{w}}_{\text{MP}}^t\|^2$  and the test error  $\|\mathbf{y}_{\text{test}} - \mathbf{X}_{\text{test}}\hat{\mathbf{w}}_{\text{MP}}^t\|^2$  as a function of iterations  $t$  (as  $t$  increases matching pursuit builds an estimate using a greater number of features). **Run the code for this part. Explain the shape of the training error plot. Does the plot for test error look similar to the one from part (i)? Comment on the similarities.** You may use conclusions obtained in previous parts to justify your comments.

### 3 CNNs on Fruits and Veggies

In this problem, we will use the dataset of fruits and vegetables that was collected in HW5. The goal is to accurately classify the produce in the image. In prior homework, we explored how to select features and then use linear classification to learn a function. We will now explore using Convolutional Neural Networks to optimize feature selection jointly with learning a classification policy.

Denote the input state  $x \in \mathbb{R}^{90 \times 90 \times 3}$ , which is a down sampled RGB image with the fruit centered in it. Each data point will have a corresponding class label, which corresponds to their matching produce. Given 25 classes, we can denote the label as  $y \in \{0, \dots, 24\}$ .

The goal of this problem is twofold. First you will learn how to implement a Convolutional Neural Network (CNN) using TensorFlow. Then we will explore some of the mysteries about why neural networks work as well as what they do in the context of a bias variance trade-off.

Note all python packages needed for the project, will be imported already. **DO NOT import new Python libraries.** Also, this project will be computationally expensive on your computer’s CPU.

Please use the free EC2 credit if you do not have a strong computer. The instructions for EC2 are on Piazza thread @391. The dataset for this project can be download [here\(link\)](#). We recommend using Tensorflow version 1.6.0 with Python3 interface.

(a) To begin the problem, we need to implement a CNN in TensorFlow. In order to reduce the burden of implementation, we are going to use a TensorFlow wrapper known as *slim*. In the starter code is a file named *cnn.py*, the network architecture and the loss function are currently blank. Using the slim library, you will have to write a convolutional neural network that has the following architecture:

- (a) Layer 1: A convolutional layer with 5 filters of size 15 by 15
- (b) Non-Linear Response: Rectified Linear Units
- (c) A max pooling operation with filter size of 3 by 3
- (d) Layer 2: A Fully Connected Layer with output size 512.
- (e) Non-Linear Response: Rectified Linear Units
- (f) Layer 3: A Fully Connected Layer with output size 25 (i.e. the class labels)
- (g) Loss Layer: Softmax Cross Entropy Loss

In the file *example\_cnn.py*, we show how to implement a network in TensorFlow Slim. Please use this as a reference. Once the network is implemented **run the script *test\_cnn\_part\_a.py* on the dataset and report the resulting confusion matrix**. The goal is to ensure that your network compiles, but we should not expect the results to be good because it is randomly initialized.

(b) The next step to train the network is to complete the pipeline which loads the datasets and offers it as mini-batches into the network. **Fill in the missing code in *data\_manager.py* and report your code**.

(c) We will now complete the iterative optimization loop. Fill in the missing code in *trainer.py* to iteratively apply SGD for a fix number of iterations. In our system, we will be using an extra momentum term to help speed up the SGD optimization. **Run the file *train\_cnn.py* and report the resulting chart**.

(d) We have seen that the convolutional neural network can achieve a good performance on this dataset. However, we are not sure about whether a fully connected neural network could do the same job. We would like to replace the layers (a)-(c) in the original network by a fully connected layer and a ReLU layer. We want to have approximately the same number of parameters between layer (a) in the original network and the new fully connected layer. **How many hidden units will the new fully connected layer has? You can round up fractional numbers to the nearest integer. Do you observe any wierd design about this new network?**

Hint: For a convolution layer with input size  $W \times H \times C$  and  $N$  convolution filters with each of them having a size of  $X \times Y$ , it has  $C \times X \times Y \times N$  parameters.



- (e) **Implement the fully connected network in part (d) and redo part (c). How's the performance of the fully connected network?**
- (f) To better understand, how the network was able to achieve the best performance on our fruits and veggies dataset. It is important to understand that it is learning features to reduce the dimensionality of the data. We can see what features were learned by examining the response maps after our convolutional layer.
- The response map is the output image after the convolutional has been applied. This image can be interpreted as what features are interesting for classification. **Fill in the missing code in `viz_features.py` and report the images specified.**
- (g) Given that our network has achieved high generalization with such low training error, it suggests that a high variance estimator is appropriate for the task. To better understand why the network is able to work, we can compare it to another high variance estimator such as K-nearest neighbors. **Fill in the missing code in `nn_classifier.py` and report the performance as the numbers of neighbors is swept across when `train_nn.py` is run.**

## 4 Running Time of $k$ -Nearest Neighbour Search Methods

The method of  $k$ -nearest neighbours is a fundamental conceptual building block of machine learning, as well as an important part of many ML algorithms. A classic example is the  $k$ -nearest neighbour classifier, which is a non-parametric classifier that finds the  $k$  closest examples in the training set to the test example and outputs the most common label among the  $k$  nearby examples as its prediction. Generating predictions using this classifier requires an algorithm to find the  $k$  closest examples in a possibly large and high-dimensional dataset, which is known as the  $k$ -nearest neighbour search problem. More precisely, given a set of  $n$  points,  $\mathcal{D} = \{\mathbf{x}_1 \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^d$  and a query point  $\mathbf{z} \in \mathbb{R}^d$ , the problem requires finding  $k$  points in  $\mathcal{D}$  that are the closest to  $\mathbf{z}$  in Euclidean distance.

$k$ -nearest neighbors is a building block that is deceptively simple — easy to understand conceptually, but tricky to analyze analytically and something that requires very serious thought to implement efficiently. This problem explores the computational complexity of nearest-neighbor methods to show how naive implementations perform very poorly as the dimensionality of the problem grows, but more sophisticated use of randomized techniques can do better.

*Overall Hint: In this problem, reading later parts will help you know what you need to do in earlier parts in case you can't figure it out. So, read ahead before asking a question.*

- (a) First, we consider the naïve exhaustive search algorithm, which exhaustively computes the distance between  $\mathbf{z}$  and all points in  $\mathcal{D}$  and returns the  $k$  points with the shortest distance. **What is the (average case) time complexity for computing distances between the query and all points, finding the  $k$  shortest distances using quickselect?** (Quickselect is a counterpart of quicksort that just picks the top  $k$  in an unordered list. Instead of taking  $O(n \log n)$  like quicksort on average, it takes  $O(n)$ . Look-up quickselect if you want, but in principle, you should be able to derive it if you understand the principle behind quicksort. Just realize that

there is no point in recursively sorting things that for sure aren't going to be in the top  $k$ .)  
**What is the (average case) time complexity of running the overall algorithm for a single query?**

- (b) Decades of research have focused on devising a way of preprocessing the data so that the  $k$ -nearest neighbours for each query can be found efficiently. “Efficient” means the time complexity of finding the  $k$ -nearest neighbours, is lower than that of the naïve exhaustive search algorithm. Since complexity of exhaustive search scales linearly in  $n$ , the complexity of a more efficient algorithm must grow more slowly than linearly as  $n$  increases; in other words, its complexity must be sublinear in  $n$ .

Many efficient algorithms for  $k$ -nearest neighbour search rely on a divide-and-conquer strategy known as space partitioning. The idea is to divide the vector space into cells and maintain a data structure that keeps track of the points that lie in each. Then, to find the  $k$ -nearest neighbours of a query, these algorithms look up the cell that contains the query and obtain the subset of points in  $\mathcal{D}$  that lie in the cell and adjacent cells. They then perform exhaustive search over this subset, i.e the distances from each point in the subset and the query are computed and the  $k$  points in the subset that are the closest to the query are returned.

For simplicity, we'll consider the special case of  $k = 1$  in the following questions, but note that the various algorithms we'll consider can be easily extended to the setting with arbitrary  $k$ . We first consider a simple partitioning scheme, where we place a Cartesian grid (a rectangular grid consisting of hypercubes) over the vector space.

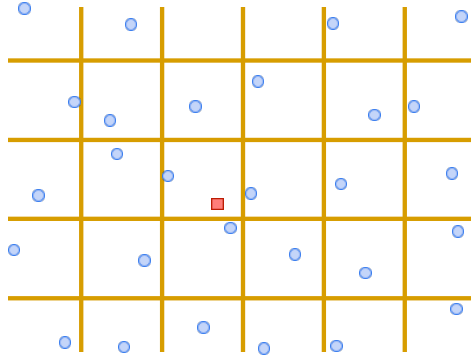


Figure 1: Illustration of the space partitioning scheme we consider. The data points are shown as blue circles and the query is shown as the red square. The cell boundaries are shown as gold lines.

**How many cells need to be searched in total if the data points are one-dimensional? Two-dimensional?  $d$ -dimensional? If each cell contains one data point, what is the time complexity for finding the 1-nearest neighbour in terms of  $d$ , assuming accessing any cell takes constant time?**

- (c) In low dimensions,  $3^d$  is much less than  $n$ , and so this method provides a significant speedup over naïve exhaustive search. However, in moderately high dimensions, because time complexity is exponential in dimensionality,  $3^d$  can easily exceed  $n$ ; in this case, the number of points retrieved from neighbouring cells is simply  $n$ , and so this method would not provide

any speedup over exhaustive search. So, if we account for the total number of points in the dataset, the query time complexity is  $O(d \min(3^d, n))$ . This exponential dependence on  $d$  arises in many settings, and is often known as *the curse of dimensionality*.

How do we overcome the curse of dimensionality? Since it arises from the need to search adjacent cells, what if we don't have cells at all?

Consider a new approach that simply projects all data points along a uniformly randomly chosen direction and keeps all projections of data points in a sorted list. To find the 1-nearest neighbour, the algorithm projects the query along the same direction used to project the data points and uses binary search to find the data point whose projection is closest to that of the query. Then it marches along the list to obtain  $\tilde{k}$  points whose projections are the closest to the projection of the query. Finally, it performs exhaustive search over these points and returns the point that is the closest to the query. This is a simplified version of an algorithm known as Dynamic Continuous Indexing (DCI).

Because this algorithm is randomized (since it uses a randomly chosen direction), there is a non-zero probability that it returns the incorrect results. We are therefore interested in how many points we need to exhaustively search over to ensure the algorithm succeeds with high probability. Equivalently, we'd like to upper bound the failure probability in terms of the number of points to search over.

The probability we consider first is the probability that a data point that is originally far away appears closer to the query under projection than a data point that is originally close. We can assume without loss of generality that the query is at the origin. Hence this probability is the same as the probability of a long vector appearing shorter than a short vector under projection. Let  $\mathbf{v}^l \in \mathbb{R}^d$  and  $\mathbf{v}^s \in \mathbb{R}^d$  denote the long and short vectors respectively and  $\mathbf{u} \in S^{d-1} \subset \mathbb{R}^d$  denote a vector drawn uniformly randomly on the unit sphere.

Assuming that  $\mathbf{0}$ ,  $\mathbf{v}^l$  and  $\mathbf{v}^s$  are not collinear, consider the plane spanned by  $\mathbf{v}^l$  and  $\mathbf{v}^s$ , which we will denote as  $P$ . (If  $\mathbf{v}^l$  and  $\mathbf{v}^s$  are collinear, random projection will essentially always be able to tell which is which so we don't bother to analyze that case. Make sure you understand why. It will help you do this problem.) For any vector  $\mathbf{w}$ , we use  $\mathbf{w}^{\parallel}$  and  $\mathbf{w}^{\perp}$  to denote the components of  $\mathbf{w}$  in  $P$  and  $P^{\perp}$  such that  $\mathbf{w} = \mathbf{w}^{\parallel} + \mathbf{w}^{\perp}$ . It should be clear that when it comes to comparing inner products, only those components that are parallel to the plane of interest matter.

For  $\mathbf{w} \in \{\mathbf{v}^s, \mathbf{v}^l\}$ , because  $\mathbf{w}^{\perp} = \mathbf{0}$ ,  $\langle \mathbf{w}, \mathbf{u} \rangle = \langle \mathbf{w}, \mathbf{u}^{\parallel} \rangle$ . So,  $\Pr\left(|\langle \mathbf{v}^l, \mathbf{u} \rangle| \leq |\langle \mathbf{v}^s, \mathbf{u} \rangle|\right) = \Pr\left(|\langle \mathbf{v}^l, \mathbf{u}^{\parallel} \rangle| \leq |\langle \mathbf{v}^s, \mathbf{u}^{\parallel} \rangle|\right)$ . **If we use  $\theta$  denote the angle of  $\mathbf{u}^{\parallel}$  relative to  $\mathbf{v}^l$ , show that this probability  $\Pr\left(|\langle \mathbf{v}^l, \mathbf{u} \rangle| \leq |\langle \mathbf{v}^s, \mathbf{u} \rangle|\right)$  is at most  $\Pr\left(|\cos \theta| \leq \|\mathbf{v}^s\|_2 / \|\mathbf{v}^l\|_2\right)$ .**

(d) **Derive the range of  $\theta$  such that  $|\cos \theta| \leq \|\mathbf{v}^s\|_2 / \|\mathbf{v}^l\|_2$  and show that**

$$\Pr\left(|\cos \theta| \leq \|\mathbf{v}^s\|_2 / \|\mathbf{v}^l\|_2\right) = 1 - \frac{2}{\pi} \cos^{-1}\left(\|\mathbf{v}^s\|_2 / \|\mathbf{v}^l\|_2\right).$$

Hint: Due to rotational invariance of a uniform distribution on the sphere, the angle between  $\mathbf{u}^{\parallel}$  and any vector in  $P$  is uniformly distributed.

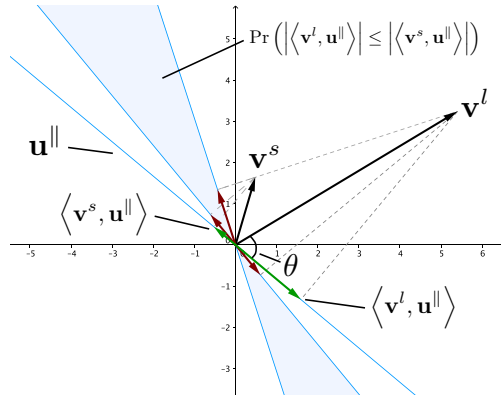


Figure 2: Examples of “good” and “bad” projection directions. The plot shows the plane  $P$  spanned by  $\mathbf{v}^l$  and  $\mathbf{v}^s$ . The vectors  $\mathbf{v}^l$  and  $\mathbf{v}^s$  are shown as black arrows. The blue lines denote possible projection directions  $\mathbf{u}^{\parallel}$ . The isolated blue line represents a “good” projection direction, since the projection of  $\mathbf{v}^l$  is longer than the projection of  $\mathbf{v}^s$  (both shown in green), thereby preserving the relative order between  $\mathbf{v}^l$  and  $\mathbf{v}^s$  in terms of their lengths after projection. Any projection direction within the shaded region is a “bad” projection direction, since the projection of  $\mathbf{v}^l$  would not be longer than the projection of  $\mathbf{v}^s$ , thereby inverting the relative order between  $\mathbf{v}^l$  and  $\mathbf{v}^s$  after projection. Two examples of such projection directions are shown, which lie on the boundaries of the shaded region. Along each of these directions, the projection  $\mathbf{v}^l$  and the projection of  $\mathbf{v}^s$  are of equal length (both of which are shown in red). The area of the shaded region represents the order-inversion probability.

- (e) The part above shows that the relative ordering of two data points is more likely to flip if their distances to the query are not very different. So, intuitively, the nearest neighbour search problem is harder if all data points are almost equidistant from the query. More concretely, we can consider a dataset consisting of many points on a sphere and a single point placed at a location we will choose later. For a query placed at the centre of the sphere, the 1-nearest neighbour search problem is easy if the single data point is placed near the centre, since it is much closer to the query than the other data points. On the other hand, if it is placed near the surface of the ball but still inside the ball, the problem is much harder since the single data point is only slightly closer to the query than the other data points and so is a much closer call. Therefore, the intrinsic hardness of the problem is characterized by the distribution of distances to the query.

The algorithm would fail to return the correct 1-nearest neighbour if more than  $\tilde{k} - 1$  points appear closer to the query than the 1-nearest neighbour under projection. The following lemma is useful:

For any set of events  $\{E_i\}_{i=1}^N$ , the probability that at least  $k'$  of them occur is at most  $\frac{1}{k'} \sum_{i=1}^N \Pr(E_i)$ .

This is a generalization of the union bound; the lemma reduces to the union bound when  $k' = 1$ . (See this paper<sup>1</sup> for the proof – let us know if you can come up with a simpler proof. We will award extra credit. :)). **Using the lemma, derive an upper bound on the probability that the algorithm fails, which is known as the *failure probability*. Use  $\mathbf{x}^{(i)}$  to denote the**

<sup>1</sup>Ke Li and Jitendra Malik. Fast  $k$ -Nearest Neighbour Search via Prioritized DCI. In *Proceedings of the 34th International Conference on Machine Learning*, pages 2081–2090, 2017.

**$i$ th closest point to the query  $\mathbf{z}$ . Then use the fact that  $1 - \frac{2}{\pi} \cos^{-1}(t) \leq t$  for all  $t \in [0, 1]$  to simplify the expression.**

Hint:

$$\Pr(\text{algorithm fails}) = \Pr\left(\text{at least } \tilde{k} \text{ points are closer to } \mathbf{z} \text{ than } \mathbf{x}^{(1)} \text{ under projection } \mathbf{u}\right)$$

- (f) Notice that the failure probability does not depend on dimensionality at all. It only depends on the distribution of distances from every point to the query, which if we recall, measures the intrinsic hardness of the problem.

What's a typical distribution of distances? Natural data usually lies on a manifold, which is a generalization of Euclidean subspace that can be "curved" (more concretely, there is a neighbourhood around every point on the manifold that resembles a low-dimensional Euclidean space). For simplicity, we'll consider the case when the data is uniformly distributed on a  $d'$ -dimensional subspace, where  $d'$  is much less than the ambient dimensionality  $d$ . Often,  $d'$  is known as the intrinsic dimensionality. Then the number of points inside a ball of radius  $r$  is roughly  $cr^{d'}$  for some constant  $c$ . So, the number of points inside a ball of constant radius grows exponentially in  $d'$ .

Assume for all  $r$  such that  $cr^{d'}$  is an integer, the number of points inside a ball centred at  $\mathbf{z}$  of radius  $r$  is exactly  $cr^{d'}$ . This is equivalent to saying  $\|\mathbf{x}^{(cr^{d'})} - \mathbf{z}\|_2 = r$  for any such  $r$ . (If we recall from the previous part,  $\mathbf{x}^{(i)}$  denotes the  $i$ th closest point to  $\mathbf{z}$ .) **Show the quantity  $\sum_{i=2}^n \|\mathbf{x}^{(1)} - \mathbf{z}\|_2 / \|\mathbf{x}^{(i)} - \mathbf{z}\|_2$  in this case is  $\sum_{i=2}^n (1/i)^{1/d'}$ .**

Hint: to derive an expression for  $\|\mathbf{x}^{(i)} - \mathbf{z}\|_2$  in terms of  $i$ , substitute  $i$  for  $cr^{d'}$  in the equality  $\|\mathbf{x}^{(cr^{d'})} - \mathbf{z}\|_2 = r$ .

- (g) (BONUS) **Show the quantity  $\sum_{i=2}^n (1/i)^{1/d'}$  is less than  $(n^{1-1/d'} - 1) / (1 - 1/d')$ .**

Hint: use the fact that  $\sum_{i=a}^b \phi(i) = \int_a^{b+1} \phi(\lfloor t \rfloor) dt$  for any function  $\phi$  and  $t - 1 < \lfloor t \rfloor$  for any  $t$ , where  $\lfloor \cdot \rfloor$  denotes the floor operator.

- (h) (BONUS) **Show the failure probability is at most  $O\left(n^{1-1/d'} / \tilde{k}\right)$  for  $d' \geq 2$ .**

- (i) If we choose the number of points to search over,  $\tilde{k}$ , to be  $\Omega\left(n^{1-1/d'}\right)$ , we can ensure that the failure probability is strictly less than 1. Choosing this value for  $\tilde{k}$  would mean that the query time complexity is  $O(dn^{1-1/d'})$ .

Consider the following variant of the algorithm, which essentially repeats the algorithm  $L$  times. More concretely, the algorithm projects the data points along  $L$  independently chosen random directions and maintain  $L$  sorted lists of projections, each corresponding to one projection direction. Given a query, we find the  $\tilde{k}$  closest points to the query along each of the  $L$  directions and exhaustively search over the union of these points.

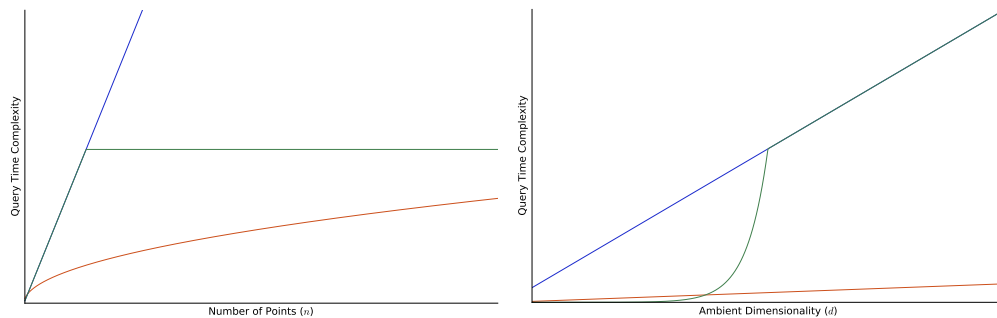
If we use  $\alpha$  denote the failure probability of the original algorithm, the failure probability of this algorithm is at most  $\alpha^L$ , because the projection directions are independently chosen. The

query time complexity is  $O(Ld(n^{1-1/d}))$ . Therefore, we can get an exponential decrease in the failure probability at a cost of a linear increase in the query time complexity.

We can choose a large enough  $L$  to make the failure probability arbitrarily small and therefore make the algorithm succeed with arbitrarily high probability. By convention,  $L$  is viewed as a constant (rather than a function of the failure probability), and so the time complexity for this algorithm is  $O(dn^{1-1/d})$ . In general, this is how the time complexity of a randomized algorithm is derived: the parameters of the algorithm are chosen so that the algorithm can succeed with arbitrarily high probability, and time complexity is computed for that choice of parameters.

Observe that the time complexity has a linear dependence on  $d$ , which is better than the exponential dependence on  $d$  of the space partitioning-based approach. The dependence on the intrinsic dimension  $d'$  is harder to see clearly since  $n$  is also involved, but notice that the penalty of making  $d'$  bigger saturates.

The following plots show the query time complexities of naïve exhaustive search, space partitioning and DCI as functions of  $n$  and  $d$ . Curves of the same colour correspond to the same algorithm. **Which algorithm does each colour correspond to?**



## 5 Your Own Question

**Write your own question, and provide a thorough solution.**

Writing your own problems is a very important way to really learn the material. The famous “Bloom’s Taxonomy” that lists the levels of learning is: Remember, Understand, Apply, Analyze, Evaluate, and Create. Using what you know to create is the top-level. We rarely ask you any HW questions about the lowest level of straight-up remembering, expecting you to be able to do that yourself. (e.g. make yourself flashcards) But we don’t want the same to be true about the highest level.

As a practical matter, having some practice at trying to create problems helps you study for exams much better than simply counting on solving existing practice problems. This is because thinking about how to create an interesting problem forces you to really look at the material from the perspective of those who are going to create the exams.

Besides, this is fun. If you want to make a boring problem, go ahead. That is your prerogative. But it is more fun to really engage with the material, discover something interesting, and then come up

with a problem that walks others down a journey that lets them share your discovery. You don't have to achieve this every week. But unless you try every week, it probably won't happen ever.